



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification⁶ :

H04L 29/06, 12/56

A1

(11) International Publication Number:

WO 99/23799

(43) International Publication Date:

14 May 1999 (14.05.99)

(21) International Application Number: PCT/GB98/03274

(22) International Filing Date: 3 November 1998 (03.11.98)

(30) Priority Data:

97308790.1

3 November 1997 (03.11.97)

EP

(71) Applicant (for all designated States except US): BRITISH TELECOMMUNICATIONS PUBLIC LIMITED COMPANY [GB/GB]; 81 Newgate Street, London EC1A 7AJ (GB).

(72) Inventors; and

(75) Inventors/Applicants (for US only): HODGKINSON, Terence, Geoffrey [GB/GB]; 46 Melton Grange Road, Melton, Woodbridge, Suffolk IP12 1SD (GB). CARTER, Simon, Francis [GB/GB]; 5 Moorfield Road, Woodbridge, Suffolk IP12 4JN (GB). O'NEILL, Alan, William [GB/GB]; 2 Rachael's Court, 36 Cemetery Road, Ipswich, Suffolk IP4 2JA (GB). WHITE, Paul, Patrick [GB/GB]; 82 Ringlow Park Road, Swinton, Manchester M27 0HB (GB).

(74) Agent: NASH, Roger, William; BT Group Legal Services, Intellectual Property Dept., Holborn Centre, 8th floor, 120 Holborn, London EC1N 2TE (GB).

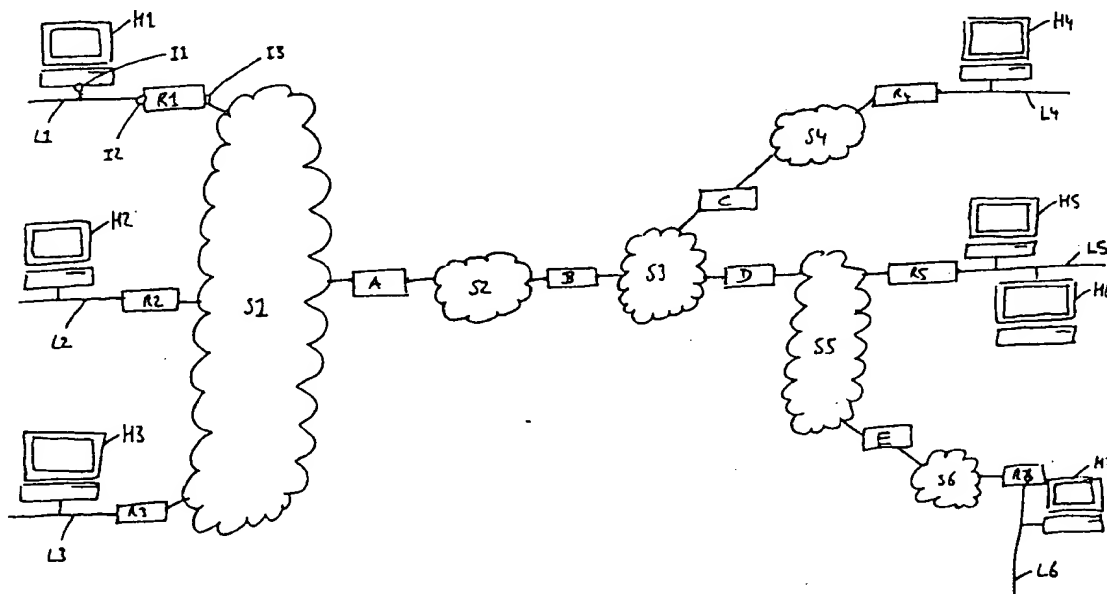
(81) Designated States: AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

Published

With international search report.

Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.

(54) Title: PACKET NETWORK



(57) Abstract

A method of reserving resources in an internet is disclosed. The method provides an improved process for use in relation to large scale shared-tree multicast environments since its use results in a reduction in path state in the routers (A-E, R1-R6) in the internet. The method involves the sending of path characteristics upstream from receivers (H1-H7) to senders (H1-H7), the routers in between combining path characteristics from different sources downstream of them. Reservations are subsequently made on the basis of the combined path characteristic data, the nature of the sender's traffic and the end-to-end quality of service required by the sender.

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

| | | | | | | | |
|----|--------------------------|----|--|----|--|----|--------------------------|
| AL | Albania | ES | Spain | LS | Lesotho | SI | Slovenia |
| AM | Armenia | FI | Finland | LT | Lithuania | SK | Slovakia |
| AT | Austria | FR | France | LU | Luxembourg | SN | Senegal |
| AU | Australia | GA | Gabon | LV | Latvia | SZ | Swaziland |
| AZ | Azerbaijan | GB | United Kingdom | MC | Monaco | TD | Chad |
| BA | Bosnia and Herzegovina | GE | Georgia | MD | Republic of Moldova | TG | Togo |
| BB | Barbados | GH | Ghana | MG | Madagascar | TJ | Tajikistan |
| BE | Belgium | GN | Guinea | MK | The former Yugoslav Republic of Macedonia | TM | Turkmenistan |
| BF | Burkina Faso | GR | Greece | | | TR | Turkey |
| BG | Bulgaria | HU | Hungary | ML | Mali | TT | Trinidad and Tobago |
| BJ | Benin | IE | Ireland | MN | Mongolia | UA | Ukraine |
| BR | Brazil | IL | Israel | MR | Mauritania | UG | Uganda |
| BY | Belarus | IS | Iceland | MW | Malawi | US | United States of America |
| CA | Canada | IT | Italy | MX | Mexico | UZ | Uzbekistan |
| CF | Central African Republic | JP | Japan | NE | Niger | VN | Viet Nam |
| CG | Congo | KE | Kenya | NL | Netherlands | YU | Yugoslavia |
| CH | Switzerland | KG | Kyrgyzstan | NO | Norway | ZW | Zimbabwe |
| CI | Côte d'Ivoire | KP | Democratic People's Republic of Korea | NZ | New Zealand | | |
| CM | Cameroon | | | PL | Poland | | |
| CN | China | KR | Republic of Korea | PT | Portugal | | |
| CU | Cuba | KZ | Kazakhstan | RO | Romania | | |
| CZ | Czech Republic | LC | Saint Lucia | RU | Russian Federation | | |
| DE | Germany | LI | Liechtenstein | SD | Sudan | | |
| DK | Denmark | LK | Sri Lanka | SE | Sweden | | |
| EE | Estonia | LR | Liberia | SG | Singapore | | |

PACKET NETWORK

The present invention relates to a method of reserving resources in a packet
5 network. It has particular utility in relation to providing an internet offering Quality
of Service guarantees.

Methods of reserving resources in an internet are well known. The method which
is currently best supported in the Internet is that defined by the Resource
10 Reservation Protocol (RSVP).

There is a desire to enable an internet to provide real-time communication as is, for
example, required if telephone conversations or video-conferences are to be
conducted over it. In this regard, two qualities of service that might be provided
15 have been specified by the following documents (which are incorporated herein by
reference):

(1) S. Schenker, C.Partridge, R.Guerin. Specification of Guaranteed Quality of
Service, Request For Comments, September 1997, RFC 2212; and
20

(2) J. Wroclawski. Specification of the Controlled-Load Network Element
Service, Request For Comments, September 1997, RFC 2211.

Essentially, Guaranteed Service as defined in the first document allows a user to
25 specify an upper bound on the time taken for his message to reach a recipient,
whereas Controlled Load Service offers a service qualitatively similar to that
provided by the internet when it is only lightly loaded. Operating an internet in
accordance with the RSVP protocol allows the provision of real-time
communication. The provision of such communication to the above-mentioned
30 Quality of Service classes when operating an internet in accordance with the RSVP
protocol is discussed in the following document (also incorporated herein by
reference):

(3) P. White. RSVP and Integrated Services in the Internet: a tutorial, IEEE Communications magazine, May 1997.

In an internetwork operating in accordance with the RSVP protocol, a sender which
5 wishes to increase the level of traffic it is sending to one or more receivers first
sends out a path information packet which contains information concerning the
characteristics of the path along which it has travelled and also information
specifying the increased traffic level. This passes through each of the nodes
through which the sender's increased traffic will pass in travelling to the one or
10 more receivers. Path characteristic data is installed at those nodes as a result of
the packet passing through them. Once this process is complete the one or more
receivers calculate a required reservation on the basis of the increased traffic
specification and end-to-end path characteristics (obtained from the received path
information packet) and send a packet back towards the sender specifying the
15 reservation required. Provided sufficient network resources are available each
node receiving the reservation-requesting packet reserves appropriate resources
and forwards the packet back towards the sender.

When an internet is operated in accordance with RSVP, receivers are responsible
20 for requesting the reservation of resources. In contrast, in an internet operating in
accordance with the Internet Stream Protocol Version 2 (ST2), senders are
responsible for requesting reservation of resources.

Although a network operating in accordance with ST2 allows resources to be
25 reserved more quickly than is possible with RSVP, it does not allow for different
levels of service to be provided in relation to a single one-to-many or many-to-
many communication.

According to the present invention there is provided a method of reserving
30 resources in a packet network comprising a plurality of hosts and one or more
interconnecting nodes, said method comprising the steps of:

operating a plurality of receiver hosts to send reverse-routed packets,
containing one or more reservation influencing parameters, along a route via one or
more of said interconnecting nodes to one or more sender hosts;

operating said one or more intermediate nodes to:
combine one or more parameters of received reverse-routed packets from
different receivers to generate combined reservation influencing parameters;
store said combined parameters; and
5 send a reverse-routed packet containing said combined parameters further
along said route towards said sender hosts;
operating said sender hosts to send a resource reservation packet to one
or more receivers back along said route; and
operating said one or more intermediate nodes, responsive to said
10 reservation packet, to reserve resources in accordance with reservation influencing
parameters stored at that node.

Because intermediate nodes examine reservation influencing data sent by
receivers, and reserve resources accordingly, the present invention better tailors a
15 resource reservation to receivers' requirements. Furthermore, the present
invention achieves this without sacrificing the faster resource reservation
associated with sender initiated resource reservation protocols.

Also, by combining path characteristic data in this way, the amount of path
20 characteristic data sent between nodes in a network where a sender sends traffic
to many receivers (i.e. where the sender is multicasting) is reduced.

In some embodiments, said one or more reservation influencing parameters include
an indication of a desired quality of service class, and said one or more
25 intermediate nodes are operable to:

update said parameters to represent the highest quality of service class
requested by downstream receivers; and

reserve resources in accordance with the resource reservation process
associated with said highest quality of service class.

30

In this way, it is possible to support different quality of service classes within a
single one-to-many or many-to-many communication. For example, some receivers
might request a quality of service in accordance with the Guaranteed Service
mentioned above, whereas others might only require a quality of service in

accordance with the Controlled Load specification mentioned above. Thus, a more flexible set of services can be provided.

In some embodiments the reservation packet carries information indicating the
5 nature of the traffic being output by the one or more senders, and node takes that information into account in calculating the resources to be reserved. In applications where the traffic from a sender is of a predetermined nature then this will not be essential.

10 In preferred embodiments, routing between the hosts is in accordance with a shared-tree protocol. In situations where there are a plurality of senders this reduces the amount of reservation influencing data that need be stored at a node.

This is especially advantageous in relation to communications involving a large
15 number of senders. A large number of senders might exist in video or audio conferences involving many people. Even greater numbers of senders are possible in multi-player games or Distributed Interactive Simulations. These potentially involve many thousands of hosts both sending and receiving messages.

20 According to a second aspect of the present invention there is provided a packet network node comprising:

means for combining one or more parameters of received reverse-routed packets to generate combined reservation influencing parameters;

means for storing said combined parameters;

25 means for sending a reverse routed packet containing said combined parameters further along said route to said sender hosts; and

means for reserving resources in accordance with combined reservation influencing parameters stored at that node responsive to a reservation packet.

30 By way of example, specific embodiments of the present invention will now be described with reference to the accompanying drawings in which:

Figure 1 is an illustration of a group of computers communicating with one another using an internet;

Figure 2 is a schematic illustration of the contents of a reservation setting packet used in a first embodiment of the present invention;

- 5 Figure 3 is a schematic illustration of the contents of an upstream-routed worst path characteristics per node packet used in the first embodiment;

Figure 4 is a schematic illustration of worst path characteristics per path data stored in a node which operates in accordance with the first embodiment; and

10

Figure 5 is a schematic illustration of worst path characteristics per interface data stored in a node which operates in accordance with the first embodiment.

- Figure 1 shows seven computers (H1 to H7), each of which is connected to a
15 Local Area Network (L1 to L6) which also includes a gateway router (R1 to R6). Two of the computers H5, H6 are connected to a single Local Area Network (LAN), the other five computers are connected to respective LANs. The gateway routers (R1 to R7) are interconnected by means of a network which comprises a number of subnets (S1 to S6) connected to one or more other subnets by
20 intermediate routers (A to E).

- A first subnet S1 connects three of the gateway routers (R1 to R3) to a first intermediate router A, which is in turn connected via a second subnet S2 to a second intermediate router B. A third subnet S3 connects the second intermediate router to third and fourth intermediate routers C, D. The third intermediate router C
25 is connected via a fourth subnet S4 to the fourth gateway router R4. The fifth subnet S5 connects the fourth intermediate router D to both the fifth intermediate router E and the fifth gateway router R5. The fifth intermediate router E is connected to the sixth gateway router R6 via a sixth subnet S6.

30

The computers (H1 to H7), routers (R1 to R7, A to E) and subnets (S1 to S6) form a portion of an internet. The internet includes a number of other computers (not shown) connected to the LANs (L1 to L6) and might include further subnets, routers and LANs. The internet is operable to allow the computers (H1 to H7) to

send messages to one another. It might, for example, be the network generally known as the Internet.

The internet operates in accordance with the TCP/IP network architecture. In addition, at least the routers and computers shown in Figure 1 support multicasting.

Each of the seven computers (H1 to H7) is controlled by a program which allows its user to participate in a multi-player game. On a user acting to influence the game, a message representing his action is sent to each of the other computers. As will be appreciated by those skilled in the art, such a message will be routed in accordance with a multicasting protocol, with the seven computers forming a multicast group. The description below assumes that a source-based tree multicast routing protocol is used, but a shared-tree protocol such as the Core Based Tree (CBT) protocol might be used instead. Those skilled in the art will have little difficulty in adapting the embodiment to operate in accordance with a shared-tree multicasting protocol.

In accordance with Internet terminology, the computers are referred to below as hosts, and the word 'node' is intended to refer to both hosts and routers. Programs running on the computers are referred to as applications. Since more than one application requiring access to the internet may be executable by any given host, each application is assigned a 'port number' so that incoming messages labelled with a port number can be passed to the corresponding application.

As those skilled in the art will know, the TCP/IP protocol suite uses addresses which refer to interfaces rather than nodes. Hence the host H1 is identified by the IP address for the interface I1 and the router R1 has one IP address for the interface via which it connects to the LAN L1 and another address for the interface I3 by which it connects to the first subnet S1.

A message is said to travel downstream from a sending host to the other hosts in the multicast group. For this reason, when discussing the reserving of resources

for a message from, say, first host H1 to the other hosts, the interface I1 between the first router R1 and the first LAN L1 is known as an upstream interface, and the interface I2 between the first router R1 and the first subnet S1 is known as a downstream interface.

5

Broadly, according to a first embodiment of the present invention, the portion of the internet serving the seven hosts (H1 to H7) is operable in accordance with the TCP/IP network architecture but is additionally operable in accordance with the following procedure to allow each of the hosts to reserve resources of the internet
10 for communication with the other hosts. By having the hosts reserve appropriate resources the characteristics of the communication paths to the other hosts can be controlled. For example, the delay in the communication can be minimised or reduced below a predetermined bound.

15 In the embodiment, each host occasionally sends a Reservation Packet (hereinafter an RES packet) to all the other hosts. Part of the purpose of this packet is to provide each node with the address of the downstream interface of the node directly upstream on the source-based tree which the message follows to each of the other hosts. This address is known as a previous hop address. Those skilled
20 in the art will recognise that a similar function is performed by the PATH messages of the RSVP protocol.

Once a recipient host has received at least one RES packet, it will start sending Backwards Control packets (hereinafter BWDC packets) towards the sending host.
25 Normally, the recipient host will send one such packet in a predetermined refresh period (though in certain circumstances (to be described below), it will send more than one).

Each of the nodes reads any BWDC packet it receives and stores worst path
30 characteristic values included in the packet in its memory. The stored values are then used to update data representing the worst per interface path characteristics (stored data like this is known as a 'state' entry to those skilled in the art). A similar worst per interface state entry is made in relation to each of any other downstream interfaces of the node. The state entries are then compared and the

worst per node path characteristics are obtained from the state entries and included in a BWDC packet to be sent upstream.

If the contents of a BWDC packet to be sent upstream differ from the contents of
5 the previous BWDC packet sent from the node, then the BWDC packet is sent immediately. Otherwise the BWDC packet is sent at the expiry of the refresh period.

This procedure results in each of the nodes in the source-based tree storing a
10 worst per interface state entry representing the worst of all the path characteristics from paths on the tree leading from the node to one or more of the recipient hosts.

On the basis of any RES packet which is subsequently received, each node
15 calculates the reservation required on each downstream link of the tree, using the corresponding worst per interface state entry.

Any installed reservations or state entries in the node are implemented using so-called 'soft-state' as is used in RSVP. This means that the state entries will be
20 deleted in the absence of appropriate refreshes known to those skilled in the art.

The reservation process will now be described in more detail in relation to Figures 2 to 5.

25 The RES packet sent from a sending host contains the information set out in Figure 2. The packet is encapsulated in an IP datagram in a similar fashion to say, a RSVP control message, being identified as a datagram relating to the present reservation protocol by a next header field (not shown) in the IP header. Those skilled in the art will be familiar with other fields in the IP header not shown in
30 Figure 2.

The other fields in the RES packet contain the following information.

- **session** - this is identical to the session field defined by the RSVP protocol - it includes the destination address (a multicast address for the multicast group) for the flow (i.e. one or more messages) to which the RES packet relates and other parameters relating to the flow. The nature of these parameters is known to those skilled in the art
- **Sender Template** - this is identical to the corresponding field defined by the RSVP protocol - i.e. it is a filter specification identifying the sending host. It contains the IP address of the sending host and optionally the sending host port being used.
- **Traffic specification (Tspec)** this is identical to the corresponding field defined by the RSVP protocol describing the sending host's traffic characteristics using the following token bucket representation.
 - p = peak rate of flow (bytes/second)
 - b = bucket depth (bytes)
 - r = token bucket rate (bytes/second)
 - m = minimum policed unit (bytes)
 - M = maximum datagram size (bytes)
- **previous hop (Phop)** - this is identical to the corresponding field defined by the RSVP protocol i.e. it is an object including the previous hop address.
- **timestamp** field - this is stamped with the time of the local node clock just before being forwarded to the next node(s) down the distribution tree.
- **end-to-end delay** field - this gives the current delay from when a packet was transmitted by the sending host until it is due to arrive at the upstream interface of the next node.
- **CRTs field(2 bits)** - this identifies the ceiling reservation type of the sending host application. 11 indicates guaranteed service, 10 indicates controlled-load, and 00 indicates best-effort. 01 is currently unspecified although may at some time

be used for a new service with quality in between best-effort and controlled-load.

- **QoSvoid bit** - if set to one this indicates that no quality of service guarantees
5 can be offered.

If the sending host requires Guaranteed Service then the datagram further contains a Guaranteed Service Object which includes the following information. The values
10 Csum and Dsum are as defined in the Guaranteed Service specification mentioned above.

- **Csum** - accumulation of C values since last upstream reshaping point.
- 15 • **Dsum** - accumulation of D values since last upstream reshaping point.
- **desired delay bound** field which indicates the end-to-end delay bound desired by the sending host application.
- 20 • **accumulated delay bound** field which indicates the installed delay bound between sending host and the upstream interface of the next node.
- **delayvoid bit** (If set, this bit is an indication to any recipient host that the desired delay bound cannot be guaranteed).
- 25 • **lossvoid bit** (If set, this bit is an indication to the recipient host that a loss-free service cannot be guaranteed).

The BWDC packets transmitted by each node and the recipient hosts running
30 under control of the application are illustrated in Figure 3. Like the RES packets, the BWDC packets are encapsulated in an IP datagram. The destination IP address is in each case the previous hop address. Those skilled in the art will realise that, in that aspect, the BWDC packets are like the RESV packets of the RSVP protocol.

The other information contained in the BWDC packet includes:

- **session** - this is identical to the session field defined above for the corresponding RES packet.
- 5
- **downstream hop object** - this is identical to RSVP next hop object - i.e. it gives the address of the upstream interface of the node directly downstream that sent the packet.
- 10
- **timestamp** field which is stamped with the time of the local node clock just before being sent to the node directly upstream.
- 15
- **timedeltaprev** field - this is filled in with the stored value of timedeltaprev (whose value is explained below) just before being sent to the node directly upstream.
- 20
- **CRT_r** field(2 bits). This field indicates the recipient host application ceiling reservation type. The mapping between the values of this field and the reservation types they represent are the same as for the CRT_s field in the RES message.
- 25
- **Worst Case Delay** field. This equals the maximum data packet propagation delay measured between the upstream interface of the node from which the BWDC packet was sent and each recipient host downstream of that node.
- 30
- **Worst Case C_{tot}** field - this is as defined in the Guaranteed Service specification - i.e. it equals the maximum accumulated C_{tot} value along the paths between the upstream interface of the node from which the BWDC packet was sent and each recipient host downstream of that node.

- **Worst Case Dtot** - this is as defined in the Guaranteed Service specification - i.e. it equals the maximum accumulated Dtot value along the paths between the upstream interface of the node from which the BWDC packet was sent and each recipient host downstream of that node.

5

- **path bandwidth** This equals the maximum path bandwidth value along the paths between the upstream interface of the node from which the BWDC packet was sent and each recipient host downstream of that node.

- 10 ◦ **sender address** - set to the address of the sending host - the address (in combination with the multicast address of the group of hosts H1 to H7) indicates to the node which source-based tree relates to the packet.

15 Ctot and Dtot are as defined in the Guaranteed Service specification. Where the sending host has requested Guaranteed Service the BWDC packet additionally includes the following:

- **excess delay field** - the amount by which the installed end-to-end delay bound currently exceeds the desired end-to-end delay bound.
- 20 ◦ **bottleneck flag** - as explained below, if set to 1 this indicates that the BWDC message has travelled at least as far as the bottleneck router (i.e. the router where the accumulated-bound first exceeded the desired-bound on the first pass of the RES message)

- 25 On setting up the reservation initially, each host transmits an RES packet to the other hosts. Each node on the source-based tree made up by the paths from the sending host to the recipient hosts carries out the following timing operations. As explained below, each node carries out these operations each time a RES packet is received.

30

Firstly, the timestamp field is read and compared to the node's local clock to determine the time it has taken the RES packet to traverse the last hop. This duration is stored at the node as the parameter `timedeltaprev`. The timestamp is

then set in accordance with the node's local clock and the RES packet is sent onwards to the next node(s) along the source-based tree towards the recipient hosts. This process is repeated until each of the nodes involved store a parameter (timedeltaprev) representing the propagation delay over the hop directly upstream
5 of them (when traversed in the downstream direction).

Once the RES packets have propagated through the nodes of the source-based tree, each of the recipient hosts sends a BWDC packet towards the sending hosts. The initial values placed in the packet by the hosts are determined as follows. The
10 path MTU, and path bandwidth correspond to the characteristics of the LAN to which the sending host is attached. The host inserts its IP address in the downstream node field and sets CRT_r in accordance with the quality of service it requires. The worst case path characteristics are all set to zero, as is the bottleneck flag. The excess delay field is unused at this stage. The timedeltaprev
15 parameter is assigned to the timedeltaprev field.

On receipt of a BWDC packet a node carries out timing operations and worst per node path characteristic determining operations as described below.

20 The timing operations involve the node first reading the timedeltaprev field in the packet. It will be realised that this records the propagation delay (experienced by the RES packet) over the hop downstream from the node (when traversed in the downstream direction). By also reading the timestamp field of the BWDC packet and the local clock a value for the propagation delay over the same hop in the
25 other direction is obtained. It is assumed that the delay over the downstream hop is the same in either direction. Hence, by taking the average of these two delays a value of the propagation delay independent of any discrepancy between the node clocks is obtained. This average delay for the downstream hop is stored by the node as a parameter 'dnext'.

30

The path characteristic monitoring operations are as follows.

Upon arrival of an BWDC packet, the node first checks for the existence of a worst per path state entry relating to the current session and to the downstream path

identified by two path identifying parameters, namely the downstream node field of the packet, and the address of the downstream interface on which the BWDC message arrived.

- 5 If no previously stored worst per path state entry is found that relates to the path and session, then a new worst per path state entry is created.

The format of a worst per path state entry is as shown in Figure 4 and includes a session identifying field, and a sender address field, the two path identifying fields
10 and additionally the following fields:

- Worst Case Ctot
- Worst Case Dtot
- Worst Case Delay
- path bandwidth
- 15 • CRTr
- pathMTU
- dnext
- bottleneck flag

- 20 The new worst per path state entry is created by assigning the CRTr, path bandwidth, pathMTU, Worst Case Delay, Worst Case Ctot and Worst Case Dtot values contained in the BWDC packet to the corresponding fields of the worst per path state entry. The parameter dnext calculated in the timing operation is assigned to the dnext field. The bottleneck flag is initialised to zero and is used as
25 explained below.

If the node has a plurality of worst per path state entries on the same interface and relate to the same session, then it combines them to create a worst per interface state entry. This situation might arise where the downstream interface connects
30 to a shared medium LAN.

As shown in Figure 5, the worst per interface state entry contains the same fields as the worst per path state entry save that it lacks the downstream node field. The values assigned to the session, sender address field and downstream interface fields are those of the worst per path state entries. The values for the other fields
5 are calculated as follows.

- Worst Case Ctot = MAX{Worst Case Ctot_i}
- Worst Case Dtot = MAX{Worst Case Dtot_i}
- Worst Case Delay = MAX{Worst Case Delay_i}
- 10 ◦ path bandwidth = MAX{path bandwidth_i}
- pathMTU = MIN{pathMTU_i}
- CRT_r = MAX{CRT_i}
- dnext = MAX{dnext_i}

15

where the subscript i takes values from 1 to the number of nodes sharing the downstream interface. It will be seen that the values Worst Case Ctot, Worst Case Dtot and Worst Case Delay that are stored in the worst per interface state entry are worst case path characteristic values for the downstream interface. It is
20 possible that the values relate to different paths from the interface. The bottleneck flag is set to one if any of the per path state entries have it set to one.

Having created a worst per interface state entry for each of its downstream interfaces, the node then generates a BWDC packet containing worst per node
25 data at least once every refresh period. The parameters to be included in the worst per node data are calculated as follows

- Worst Case Ctot = MAX{Worst Case Ctot_n + Clocal_n} where Clocal_n is the value of Ctot between the downstream interface on which the BWDC packet arrived
30 and the upstream interface (determined by the session) of the node;

- Worst Case $D_{tot} = \text{MAX}\{\text{Worst Case } D_{tot_n} + D_{local_n}\}$ where D_{local_n} is the value of C_{tot} between the downstream interface on which the BWDC packet arrived and the upstream interface (determined by the session) of the node;
- 5 • Worst Case Delay = $\text{MAX}\{\text{Worst Case Delay}_n\} + d_{next_n}$
- path bandwidth = $\text{MIN}(\text{MAX}\{\text{path bandwidth}_n\}, \text{path bandwidth upstream})$
 - $CRT_r = \text{MAX}\{CRT_n\}$
- 10 • pathMTU = $\text{MIN}(\text{MIN}\{\text{pathMTU}_n\}, \text{pathMTU upstream})$

where n is an index which takes values from one to the number of downstream interfaces from the node and 'pathMTU upstream' and 'path bandwidth upstream' represent the minimum packet size and the bandwidth of the upstream hop respectively.

15

If CRT_r is not set to 11 then C_{tot} , D_{tot} and Worst Case Delay are set to zero.

- 20 The above operations are repeated by each node in the source-based tree relating to the current session. It will be realised that for a flow from a given sending host, each node stores the worst case characteristics of all the paths leading from the node to the recipient hosts.
- 25 The sending host then sends a further RES packet (Figure 2) requesting, for example, Guaranteed Service. Hence, the CRT_r field in the packet is set to 11 and the desired delay bound is set to the limit that the receiver wishes to be placed on the time taken for a packet from the sending host to reach all the recipient hosts.
- 30 The Traffic Specification fields are filled in by the sending host, and the end-to-end delay field is set to the value of the parameter d_{next} stored at the sending host. Each of the fields of the Guaranteed Service other than the desired delay bound is initialised to zero.

At the first node downstream along the source-based tree, the queuing delay to be imposed on the flow in relation to each of the one or more downstream hops involved in the current session. This is determined in accordance with the equation below:

5

Equation (1)

$Q_{delay} = \text{desired-bound} - \text{accumulated-bound} - \text{Worst Case Delay}$

- 10 The Worst Case Delay value equals the sum of the corresponding field and the dnext parameter of the worst per interface state entry at the current node. At the first node the accumulated bound will normally be zero or negligible. The Q_{delay} value represents an estimate of the total queuing delay that can be tolerated over the remainder of the path to a recipient host.

15

Processing similar to that carried out by receivers operating in accordance with the RSVP protocol is then used to calculate a bandwidth to be reserved on each of the downstream hops in order to stay 'on-course' for the desired bound. This calculation will often be an overestimate because it may be that the individual
 20 worst per interface parameters relate to different paths from that interface. Those skilled in the art will realise that the estimate is calculated using the following equations:

Equation 2

$$25 \quad Q_{delay_{end2end}} = \frac{(b - M)(p - R)}{R(p - r)} + \frac{(M + C_{tot})}{R} + D_{tot} \quad (\text{case } p > R \geq r)$$

Equation 3

$$Q_{delay_{end2end}} = \frac{(M + C_{tot})}{R} + D_{tot} \quad (\text{case } R \geq p \geq r)$$

- 30 The parameters M, p, b and r are obtained from the corresponding fields of the RES message (the meaning of those symbols is as set out above in relation to the RES message). The values of C_{tot} and D_{tot} are a sum of the values from the

worst per interface state entry relating to the current downstream interface and the local values of C and D. To obtain the value of Qdelay to insert into the equations 2 and 3, equation 1 is used. The equation is solved to yield a value for R and a bandwidth R is reserved over the hop leading from the current downstream
5 interface.

In accordance with the rules set out in the Guaranteed Service specification, the values of CSum and DSum are used along with other parameters to calculate the amount of buffering that must be assigned to the reservation to prevent packet
10 loss.

Once the reservation request has been serviced the end-to-end delay field in the RES packet is increased by adding to it the following:

- 15
- The propagation delay, dnext for the next hop
 - An estimate of the current local queuing delay (for the relevant outgoing interface) for data packets of the flow to which the RES packet refers.

The following updates are also made to the RES packet:

20

- The accumulated-bound field of the copied packet is increased by:

1) The propagation delay, dnext for the next hop, and

2) The installed local queuing delay bound (for the relevant downstream
25 interface) for messages of the flow to which the RES packet refers. This local queuing delay bound is obtained by inserting the reserved value of R into Equation 2 and Equation 3 along with the local values C and D for the parameters Ctot and Dtot respectively.

30 The following updates are then made to the RES packet

- the local value of C is added to the CSum field
- the local value of D is added to the DSum field

unless re-shaping of the sending hosts traffic is carried out at the downstream interface, in which case both CSum and Dsum are set to zero.

Also, the Qosvoid field is set to one if either

- 5 a) a Guaranteed Service reservation attempt fails to reserve a bandwidth at least as great as the token bucket rate of the traffic specification; or
- b) a Controlled Load reservation attempt fails.

If a node receives a RES packet with Qosvoid set to one, then if $\text{MIN}\{\text{CRT}_s, \text{CRT}_r\}$ is 10 or 11 it attempts to secure a Controlled Load reservation.

10

Once updating of the fields is complete the timestamp field is (as described above) set equal to the local clock before forwarding the RES packet to each next hop down the routing tree.

- 15 Similar processing is carried out at each of the subsequent nodes, the increase of the accumulated bound by the installed queuing delay resulting in the value of Qdelay obtained from equation (1) continuing to be an estimate of the total queuing delay that can be tolerated for the remainder of the path to a recipient host. Assuming the messages from the host to be in accordance with the declared
- 20 traffic specification, and provided each of the nodes is able to reserve the calculated bandwidth R, the above-described processing at the nodes will be effective to enable messages from the sending hosts to be delivered within the desired delay bound to all of the recipient hosts.

- 25 If any node is unable to reserve the bandwidth as calculated above then (subject to policy configurations at the node) as much bandwidth as possible is reserved. If, however, the amount of bandwidth reserved is less than the value of the token bucket rate field of the RES packet then Guaranteed Service is not offered and this is indicated by setting the delayvoid, lossvoid and Qosvoid flags to 1.

30

If the bandwidth reserved is more than the value of the token bucket rate field of the RES packet, then the accumulated bound is compared to the desired bound.

Where the desired bound is higher then the reservation process continues as above. Since the reserved bandwidth R may be overestimated as explained above, this often results in the setting up of reservations which provide a delay over the end-to-end path which is less than the desired delay bound .

5

Where, however, the desired bound is lower, the node

- sets a reservation bottleneck flag associated with the relevant per interface state entry, thereby labelling itself as a bottleneck node.

10

- sets the delayvoid flag in the RES packet to one and forwards it to nodes further down the tree.

Where a node receives an RES packet with the delayvoid field set to one, it
15 reserves the maximum bandwidth possible (perhaps limited to a multiple of the peak rate of the sending host) and forwards it on.

On receiving an RES packet with delayvoid field set to 1 and delayloss field set to 0 (for $CRTs = CRT_r = 11$) the recipient host calculates the amount by which the
20 accumulated delay bound field exceeds the desired delay bound field. The recipient host then immediately sends a BWDC packet with the excess delay field set to the calculated excess delay and with the bottleneck field set to zero. Each node receiving such an BWDC packet ignores the packet unless the bottleneck flag of the node's associated worst per interface state entry is set to 1. The bottleneck
25 node receives the BWDC packet and sets the bottleneck field to one. The node attempts to eliminate or reduce the excess delay indicated in the BWDC message by increasing the local bandwidth reservation. Following this if excess delay still exists the BWDC packet with the modified value of excess delay is then sent a hop at a time towards the sender with a reservation increase being attempted at each
30 node until either the BWDC packet reaches the sender or the excess delay is eliminated.

In the above embodiment, data is combined at a node and sent to a node upstream where that combined data is used to calculate an appropriate resource reservation.

In other embodiments the resource reservation might be calculated on the basis of combined data at the node where the combination takes place.

A second embodiment is similar to the first embodiment but uses a shared tree
5 protocol. In the second embodiment the sender address field (in the BWDC
packet, the worst per path and worst per interface state entries) is not required.
This is because the interfaces out of which the RES packet are to be sent are
determinable by the node on the basis of the session parameter and the incoming
interface. Since, when using a shared tree the worst per interface state entries
10 will relate to the multicast group rather than a specific sender to the multicast
group, the number of worst per interface state entries in each node is only as great
as the number of outgoing interfaces from the node.

It will be seen how, in multi-sender environments, the second embodiment reduces
15 the amount of worst per interface state entries required in comparison with known
reservation protocols.

In the above embodiments, the worst-case merging of C terms (e.g. C_{tot}), D terms
(e.g. D_{tot}) and link propagation delay was carried out for each term independently.
20 This results in an overly conservative local bandwidth reservation. In preferred
embodiments, a rate independent delay parameter which includes both the D term
and the link propagation delay is used. This might be done by taking the
forwarded value of D as the value from the worst case rate independent delay
parameter rather than simply the maximum D value from each path.

25

It will be seen how the above embodiments enable a router to find the highest
quality of service requested by any one of the downstream receivers. Similar
considerations apply to path characteristic data and other reservation influencing
parameters. It is the combination of such parameters at intermediate nodes in the
30 network that allows a more flexible resource reservation to be provided for
multicast communication in a packet network.

CLAIMS

1. A method of reserving resources in a packet network comprising a plurality of hosts and one or more interconnecting nodes, said method comprising
5 the steps of:
 - operating a plurality of receiver hosts to send reverse-routed packets, containing one or more reservation influencing parameters, along a route via one or more of said interconnecting nodes to one or more sender hosts;
 - operating said one or more intermediate nodes to:
10 combine one or more parameters of received reverse-routed packets from different receivers to generate combined reservation influencing parameters;
 - store said combined parameters; and
 - send a reverse-routed packet containing said combined parameters further
15 along said route towards said sender hosts;
 - operating said sender hosts to send a resource reservation packet to one or more receivers back along said route; and
 - operating said one or more intermediate nodes, responsive to said reservation packet, to reserve resources in accordance with reservation influencing parameters stored at that node.
20
2. A method according to claim 1 wherein:
 - said one or more reservation influencing parameters include an indication of a desired quality of service class;
 - said one or more intermediate nodes are operable to:
25 update said parameters to represent the highest quality of service class requested by downstream receivers; and
 - reserve resources in accordance with the resource reservation process associated with said highest quality of service class.
- 30 3. A method according to claim 1 wherein said one or more reservation influencing parameters include path characteristic parameters.
4. A method according to claim 3 wherein said path characteristic parameters comprise one or more delay parameters.

5. A method according to any preceding claim wherein said hosts communicate in accordance with a shared tree routing algorithm.
- 5 6. A packet network node comprising:
means for combining one or more parameters of received reverse-routed packets to generate combined reservation influencing parameters;
means for storing said combined parameters;
means for sending a reverse routed packet containing said combined
10 parameters further along said route to said sender hosts; and
means for reserving resources in accordance with combined reservation influencing parameters stored at that node responsive to a reservation packet.
7. A method of reserving resources in a packet network comprising a
15 plurality of hosts and one or more interconnecting nodes, said method comprising the steps of:
operating one or more receiver hosts to send path characteristic data packets, containing one or more path parameters, along a route via one or more of said interconnecting nodes to one or more sender hosts;
20 operating said one or more intermediate nodes to:
process one or more parameters of received path characteristic data packets to generate updated path characteristic parameters;
store said updated parameters; and
send a path characteristic data packet containing said updated parameters
25 further along said route to said sender hosts
operating said sender hosts to send a resource reservation packet to one or more receivers back along said route; and
operating said one or more intermediate nodes, responsive to said reservation packet, to reserve resources in accordance with path characteristic
30 data stored at that node.

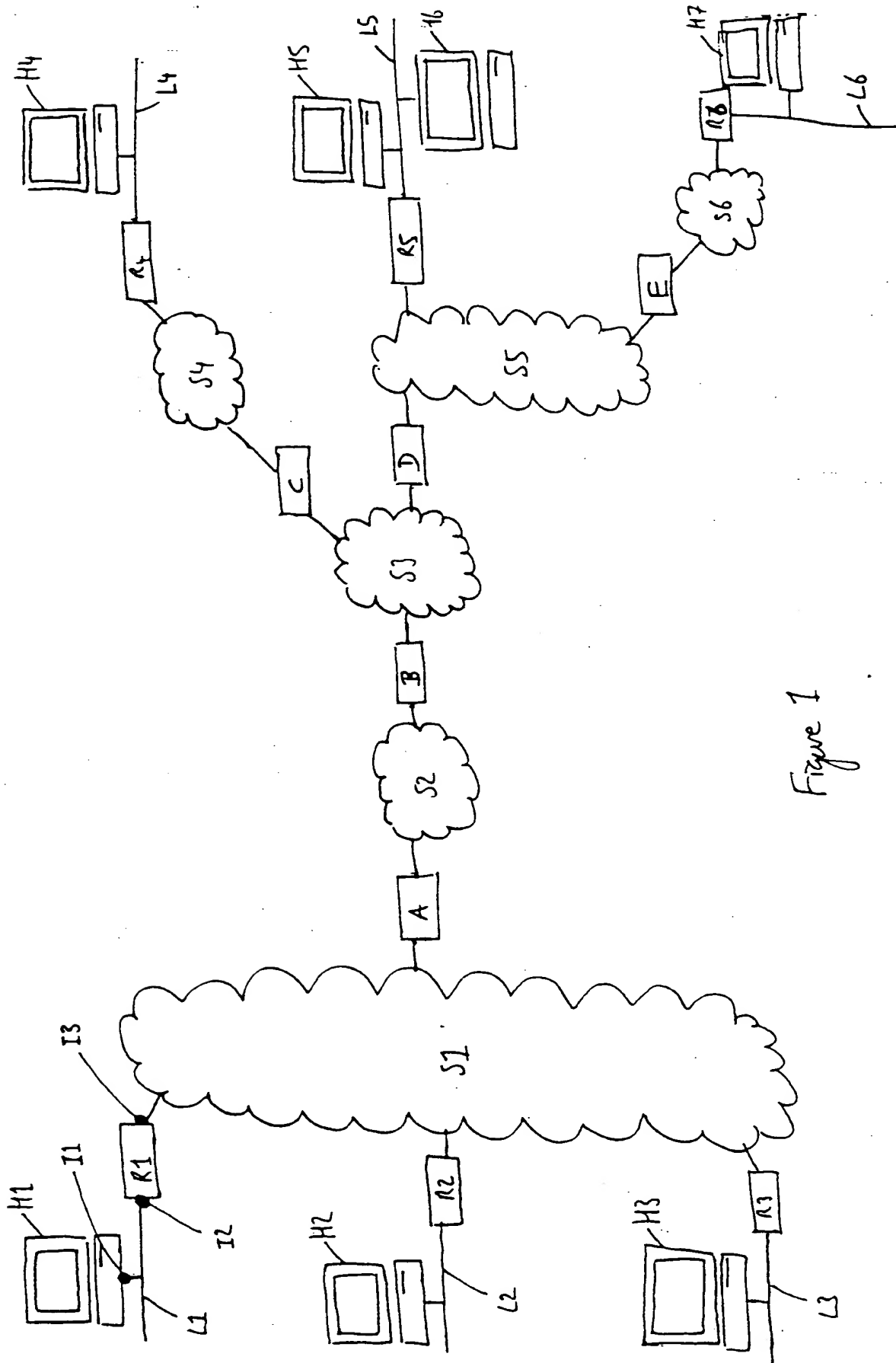


Figure 1

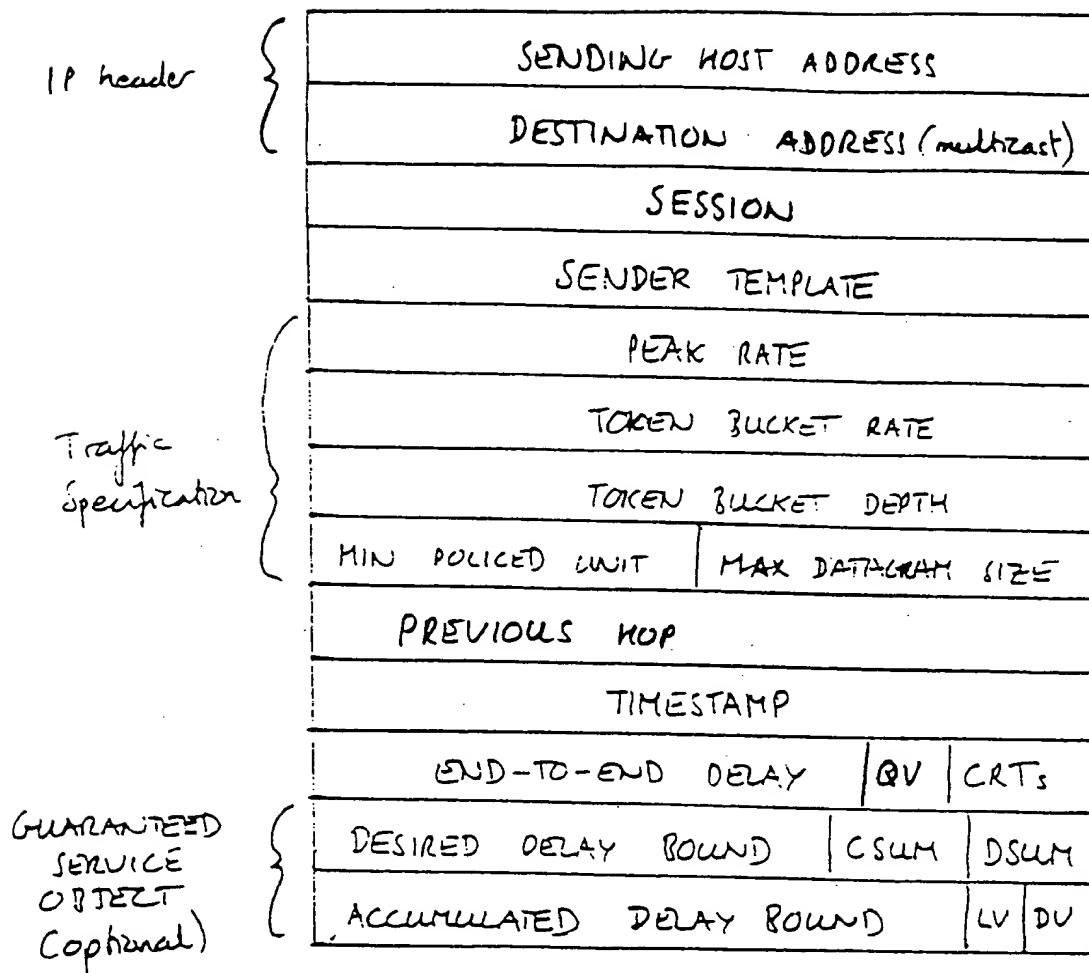


Figure 2

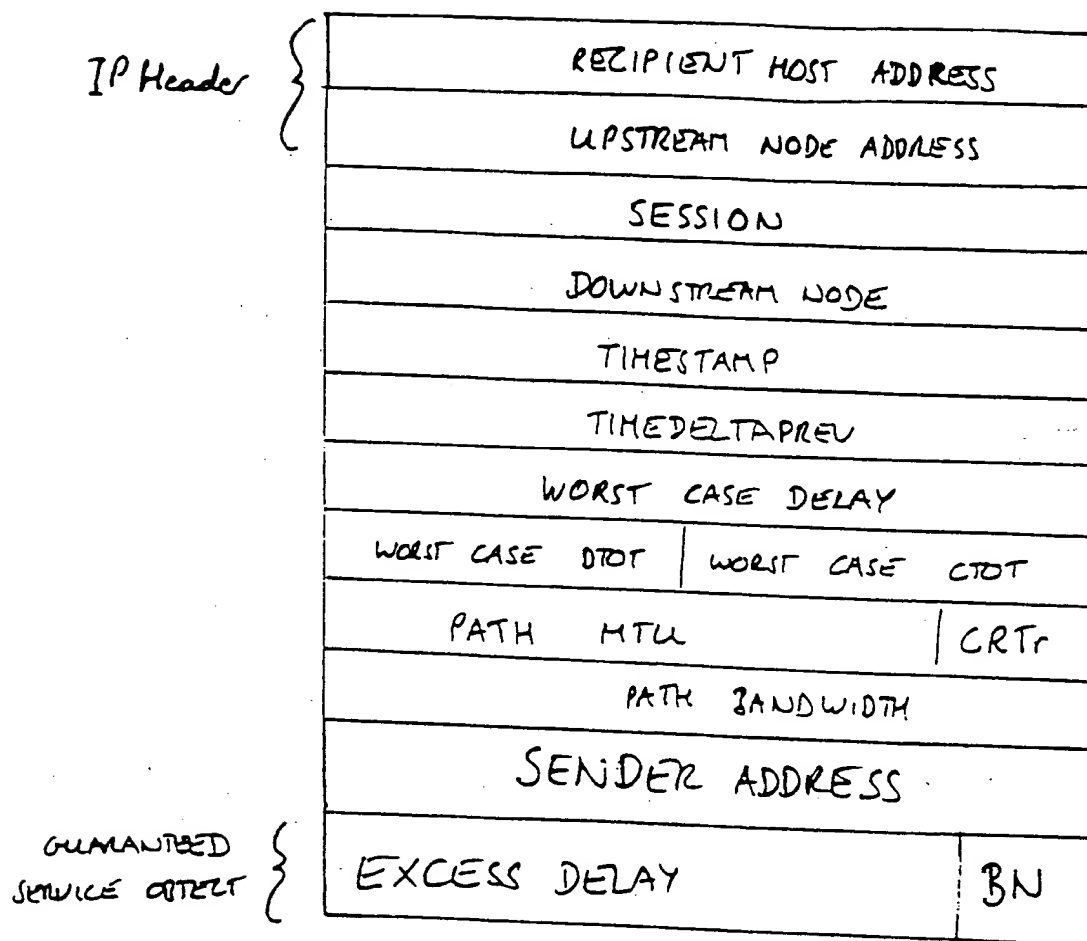


Figure 3

| | |
|----------------------|------------------|
| SESSION | |
| DOWNSTREAM NODE | |
| DOWNSTREAM INTERFACE | |
| WORST CASE CDT | WORST CASE DTOT |
| WORST CASE DELAY | |
| PATH | BANDWIDTH |
| PATH MTU | CRT _r |
| DNEXT | BN |
| SENDER ADDRESS | |

Figure 4

| | |
|----------------------|------------------|
| SESSION | |
| DOWNSTREAM INTERFACE | |
| WORST CASE CTOT | WORST CASE DTOT |
| WORST CASE DELAY | |
| PATH BANDWIDTH | |
| PATH MTU | CRT _r |
| DNEXT | BN |
| SENDER ADDRESS | |

Figure 5

INTERNATIONAL SEARCH REPORT

Inte. onal Application No

PCT/GB 98/03274

A. CLASSIFICATION OF SUBJECT MATTER
IPC 6 H04L29/06 H04L12/56

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 6 H04L

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|------------|---|-----------------------|
| X | SHENKER S ET AL: "TWO ISSUES IN RESERVATION ESTABLISHMENT" COMPUTER COMMUNICATIONS REVIEW, vol. 25, no. 4, 1 October 1995, pages 14-26, XP000541647 | 1-6 |
| A | see paragraph 2.1 - paragraph 2.3 see paragraph 3.1.2 see paragraph 3.3 --- -/-- | 7 |

☒ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

Date of the actual completion of the international search

12 March 1999

Date of mailing of the international search report

26/03/1999

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Dupuis, H

INTERNATIONAL SEARCH REPORT

Int. onal Application No

PCT/GB 98/03274

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

| Category | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|----------|---|-----------------------|
| X | <p>DELGROSSI L AND BERGER L: "RFC 1819" INTERNET STREAM PROTOCOL VERSION 2 (ST2) PROTOCOL SPECIFICATION, August 1995, pages 28-45, XP002060637 INTERNET ENGINEERING TASK FORCE see paragraph 4.1.1 see paragraph 4.2 see paragraph 4.5.2 - paragraph 4.5.7 see paragraph 4.5.10 see paragraph 4.6.1 - paragraph 4.6.3 -----</p> | 7 |